

# Python For Data Science

## python™ Basics Cheat Sheet

Learn Python Basics online at [www.DataCamp.com](http://www.DataCamp.com)

## > Variables and Data Types

### Variable Assignment

```
>>> x=5
>>> x
5
```

### Calculations With Variables

```
>>> x+2 #Sum of two variables
7
>>> x-2 #Subtraction of two variables
3
>>> x*2 #Multiplication of two variables
10
>>> x**2 #Exponentiation of a variable
25
>>> x%2 #Remainder of a variable
1
>>> x/float(2) #Division of a variable
2.5
```

### Types and Type Conversion

```
str()
'5', '3.45', 'True' #Variables to strings

int()
5, 3, 1 #Variables to integers

float()
5.0, 1.0 #Variables to floats

bool()
True, True, True #Variables to booleans
```

## > Libraries

pandas	NumPy	matplotlib	learn
Data analysis	Scientific computing	2D plotting	Machine learning

### Import Libraries

```
>>> import numpy
>>> import numpy as np
```

### Selective import

```
>>> from math import pi
```

## > Strings

```
>>> my_string = 'thisStringIsAwesome'
>>> my_string
'thisStringIsAwesome'
```

### String Operations

```
>>> my_string * 2
'thisStringIsAwesomethisStringIsAwesome'
>>> my_string + 'Innit'
'thisStringIsAwesomeInnit'
>>> 'm' in my_string
True
```

### String Indexing

Index starts at 0

```
>>> my_string[3]
>>> my_string[4:9]
```

### String Methods

```
>>> my_string.upper() #String to uppercase
>>> my_string.lower() #String to lowercase
>>> my_string.count('w') #Count String elements
>>> my_string.replace('e', 'i') #Replace String elements
>>> my_string.strip() #Strip whitespaces
```

## > NumPy Arrays

Also see Lists

```
>>> my_list = [1, 2, 3, 4]
>>> my_array = np.array(my_list)
>>> my_2darray = np.array([[1,2,3],[4,5,6]])
```

### Selecting Numpy Array Elements

Index starts at 0

```
Subset
>>> my_array[1] #Select item at index 1
2

Slice
>>> my_array[0:2] #Select items at index 0 and 1
array([1, 2])

Subset 2D Numpy arrays
>>> my_2darray[:,0] #my_2darray[rows, columns]
array([1, 4])
```

### Numpy Array Operations

```
>>> my_array > 3
array([False, False, False, True], dtype=bool)
>>> my_array * 2
array([2, 4, 6, 8])
>>> my_array + np.array([5, 6, 7, 8])
array([6, 8, 10, 12])
```

### Numpy Array Functions

```
>>> my_array.shape #Get the dimensions of the array
>>> np.append(other_array) #Append items to an array
>>> np.insert(my_array, 1, 5) #Insert items in an array
>>> np.delete(my_array,[1]) #Delete items in an array
>>> np.mean(my_array) #Mean of the array
>>> np.median(my_array) #Median of the array
>>> my_array.corrcoef() #Correlation coefficient
>>> np.std(my_array) #Standard deviation
```

## > Lists

Also see NumPy Arrays

```
>>> a = 'is'
>>> b = 'nice'
>>> my_list = ['my', 'list', a, b]
>>> my_list2 = [[4,5,6,7], [3,4,5,6]]
```

### Selecting List Elements

Index starts at 0

Subset

```
>>> my_list[1] #Select item at index 1
>>> my_list[-3] #Select 3rd last item
```

Slice

```
>>> my_list[1:3] #Select items at index 1 and 2
>>> my_list[1:] #Select items after index 0
>>> my_list[:3] #Select items before index 3
>>> my_list[:] #Copy my_list
```

Subset Lists of Lists

```
>>> my_list2[1][0] #my_list[list][itemOfList]
>>> my_list2[1][:2]
```

### List Operations

```
>>> my_list + my_list
['my', 'list', 'is', 'nice', 'my', 'list', 'is', 'nice']
>>> my_list * 2
['my', 'list', 'is', 'nice', 'my', 'list', 'is', 'nice']
>>> my_list2 > 4
True
```

### List Methods

```
>>> my_list.index(a) #Get the index of an item
>>> my_list.count(a) #Count an item
>>> my_list.append('!) #Append an item at a time
>>> my_list.remove('!) #Remove an item
>>> del(my_list[0:1]) #Remove an item
>>> my_list.reverse() #Reverse the list
>>> my_list.extend('!) #Append an item
>>> my_list.pop(-1) #Remove an item
>>> my_list.insert(0,'!) #Insert an item
>>> my_list.sort() #Sort the list
```

## > Python IDEs (Integrated Development Environment)

ANACONDA.

Leading open data science platform powered by Python

SPYDER

Free IDE that is included with Anaconda

jupyter

Create and share documents with live code

## > Asking For Help

```
>>> help(str)
```

Learn Data Skills Online at  
[www.DataCamp.com](http://www.DataCamp.com)

# Python For Data Science

## NumPy Cheat Sheet

Learn NumPy online at [www.DataCamp.com](http://www.DataCamp.com)

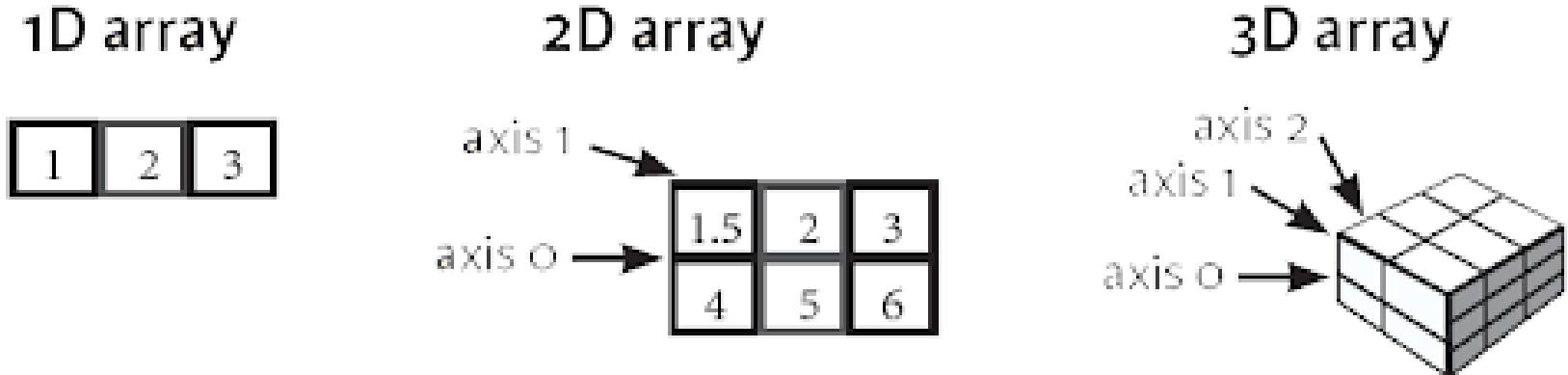
## Numpy

The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.

Use the following import convention:

```
>>> import numpy as np
```

## NumPy Arrays



## Creating Arrays

```
>>> a = np.array([1,2,3])
>>> b = np.array([(1.5,2,3), (4,5,6)], dtype = float)
>>> c = np.array([(1.5,2,3), (4,5,6)],[(3,2,1), (4,5,6)]), dtype = float)
```

## Initial Placeholders

```
>>> np.zeros((3,4)) #Create an array of zeros
>>> np.ones((2,3,4),dtype=np.int16) #Create an array of ones
>>> d = np.arange(10,25,5) #Create an array of evenly spaced values (step value)
>>> np.linspace(0,2,9) #Create an array of evenly spaced values (number of samples)
>>> e = np.full((2,2),7) #Create a constant array
>>> f = np.eye(2) #Create a 2x2 identity matrix
>>> np.random.random((2,2)) #Create an array with random values
>>> np.empty((3,2)) #Create an empty array
```

## I/O

### Saving & Loading On Disk

```
>>> np.save('my_array', a)
>>> np.savez('array.npz', a, b)
>>> np.load('my_array.npy')
```

### Saving & Loading Text Files

```
>>> np.loadtxt("myfile.txt")
>>> np.genfromtxt("my_file.csv", delimiter=',')
>>> np.savetxt("myarray.txt", a, delimiter=" ")
```

## Asking For Help

```
>>> np.info(np.ndarray.dtype)
```

## Inspecting Your Array

```
>>> a.shape #Array dimensions
>>> len(a) #Length of array
>>> b.ndim #Number of array dimensions
>>> e.size #Number of array elements
>>> b.dtype #Data type of array elements
>>> b.dtype.name #Name of data type
>>> b.astype(int) #Convert an array to a different type
```

## Data Types

```
>>> np.int64 #Signed 64-bit integer types
>>> np.float32 #Standard double-precision floating point
>>> np.complex #Complex numbers represented by 128 floats
>>> np.bool #Boolean type storing TRUE and FALSE values
>>> np.object #Python object type
>>> np.string_ #Fixed-length string type
>>> np_unicode_ #Fixed-length unicode type
```

## Array Mathematics

### Arithmetic Operations

```
>>> g = a - b #Subtraction
array([-0.5, 0. , 0. ],
      [-3. , -3. , -3. ])
>>> np.subtract(a,b) #Subtraction
>>> b + a #Addition
array([[ 2.5, 4. , 6. ],
      [ 5. , 7. , 9. ]])
>>> np.add(b,a) Addition
>>> a / b #Division
array([[ 0.66666667, 1. , 1. ],
      [ 0.25 , 0.4 , 0.5 ]])
>>> np.divide(a,b) #Division
>>> a * b #Multiplication
array([[ 1.5, 4. , 9. ],
      [ 4. , 10. , 18. ]])
>>> np.multiply(a,b) #Multiplication
>>> np.exp(b) #Exponentiation
>>> np.sqrt(b) #Square root
>>> np.sin(a) #Print sines of an array
>>> np.cos(b) #Element-wise cosine
>>> np.log(a) #Element-wise natural logarithm
>>> e.dot(f) #Dot product
array([[ 7., 7.],
      [ 7., 7.]])
```

### Comparison

```
>>> a == b #Element-wise comparison
array([[False, True, True],
      [False, False, False]], dtype=bool)
>>> a < 2 #Element-wise comparison
array[[True, False, False], dtype=bool]
>>> np.array_equal(a, b) #Array-wise comparison
```

### Aggregate Functions

```
>>> a.sum() #Array-wise sum
>>> a.min() #Array-wise minimum value
>>> b.max(axis=0) #Maximum value of an array row
>>> b.cumsum(axis=1) #Cumulative sum of the elements
>>> a.mean() #Mean
>>> np.median(b) #Median
>>> np.correlcoef(a) #Correlation coefficient
>>> np.std(b) #Standard deviation
```

## Copying Arrays

```
>>> h = a.view() #Create a view of the array with the same data
>>> np.copy(a) #Create a copy of the array
>>> h = a.copy() #Create a deep copy of the array
```

## Sorting Arrays

```
>>> a.sort() #Sort an array
>>> c.sort(axis=0) #Sort the elements of an array's axis
```

## Subsetting, Slicing, Indexing

### Subsetting

```
>>> a[2] #Select the element at the 2nd index
3
>>> b[1,2] #Select the element at row 1 column 2 (equivalent to b[1][2])
6.0
```

1	2	3
1.5	2	3
4	5	6

### Slicing

```
>>> a[0:2] #Select items at index 0 and 1
array([1, 2])
>>> b[0:2,1] #Select items at rows 0 and 1 in column 1
array([ 2., 2., 3.])
>>> b[:,1] #Select all items at row 0 (equivalent to b[0:, 1])
array([[1.5, 2., 3., 1.5]])
>>> c[1,...] #Same as [1,:,:]
array([[ 3., 2., 1.,
        4., 5., 6.]]))
>>> a[ : :-1] #Reversed array a array([3, 2, 1])
```

1	2	3
1.5	2	3
4	5	6

### Boolean Indexing

```
>>> a[a<2] #Select elements from a less than 2
array([1])
```

### Fancy Indexing

```
>>> b[[1, 0, 1, 0],[0, 1, 2, 0]] #Select elements (1,0),(0,1),(1,2) and (0,0)
array([ 4., 2., 6., 1.5])
>>> b[[1, 0, 1, 0]][:, [0,1,2,0]] #Select a subset of the matrix's rows and columns
array([[ 4., 5., 6., 4.],
      [ 1.5, 2., 3., 1.5],
      [ 4., 5., 6., 4.],
      [ 1.5, 2., 3., 1.5]])
```

1	2	3
1.5	2	3
4	5	6

## Array Manipulation

### Transposing Array

```
>>> i = np.transpose(b) #Permute array dimensions
>>> i.T #Permute array dimensions
```

### Changing Array Shape

```
>>> b.ravel() #Flatten the array
>>> g.reshape(3,-2) #Reshape, but don't change data
```

### Adding/Removing Elements

```
>>> h.resize((2,6)) #Return a new array with shape (2,6)
>>> np.append(h,g) #Append items to an array
>>> np.insert(a, 1, 5) #Insert items in an array
>>> np.delete(a,[1]) #Delete items from an array
```

### Combining Arrays

```
>>> np.concatenate((a,d),axis=0) #Concatenate arrays
array([[ 1, 2, 3, 10],
      [ 2, 15],
      [ 3, 20]])
>>> np.vstack((a,b)) #Stack arrays vertically (row-wise)
array([[ 1, 2, 3, 10],
      [ 1.5, 2, 3, 1.5],
      [ 4., 5., 6., 4.],
      [ 1.5, 2., 3., 1.5]])
>>> np.r_[e,f] #Stack arrays vertically (row-wise)
>>> np.hstack((e,f)) #Stack arrays horizontally (column-wise)
array([[ 7., 7., 1., 0.],
      [ 7., 7., 0., 1.]])
>>> np.column_stack((a,d)) #Create stacked column-wise arrays
array([[ 1, 10],
      [ 2, 15],
      [ 3, 20]])
>>> np.c_[a,d] #Create stacked column-wise arrays
```

### Splitting Arrays

```
>>> np.hsplit(a,3) #Split the array horizontally at the 3rd index
[array([1]),array([2]),array([3])]
>>> np.vsplit(c,2) #Split the array vertically at the 2nd index
[array([[ 1.5, 2., 1.],
        [ 4., 5., 6.]]),
 array([[ 3., 2., 3.],
        [ 4., 5., 6.]]])
```

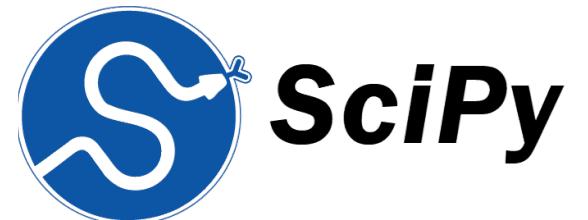


# Python For Data Science

## SciPy Cheat Sheet

Learn SciPy online at [www.DataCamp.com](http://www.DataCamp.com)

### SciPy



The SciPy library is one of the core packages for scientific computing that provides mathematical algorithms and convenience functions built on the NumPy extension of Python.

### > Interacting With NumPy

Also see NumPy

```
>>> import numpy as np
>>> a = np.array([1,2,3])
>>> b = np.array([(1+5j),2j,3j], [(4j,5j,6j)])
>>> c = np.array([[1.5,2,3], [4,5,6]], [(3,2,1), (4,5,6)])
```

#### Index Tricks

```
>>> np.mgrid[0:5,0:5] #Create a dense meshgrid
>>> np.ogrid[0:2,0:2] #Create an open meshgrid
>>> np.r_[[3,[0]*5,-1:1:10j]] #Stack arrays vertically (row-wise)
>>> np.c_[b,c] #Create stacked column-wise arrays
```

#### Shape Manipulation

```
>>> np.transpose(b) #Permute array dimensions
>>> b.flatten() #Flatten the array
>>> np.hstack((b,c)) #Stack arrays horizontally (column-wise)
>>> np.vstack((a,b)) #Stack arrays vertically (row-wise)
>>> np.hsplit(c,2) #Split the array horizontally at the 2nd index
>>> np.vsplit(d,2) #Split the array vertically at the 2nd index
```

#### Polynomials

```
>>> from numpy import poly1d
>>> p = poly1d([3,4,5]) #Create a polynomial object
```

#### Vectorizing Functions

```
>>> def myfunc(a):
...     if a < 0:
...         return a*2
...     else:
...         return a/2
>>> np.vectorize(myfunc) #Vectorize functions
```

#### Type Handling

```
>>> np.real(c) #Return the real part of the array elements
>>> np.imag(c) #Return the imaginary part of the array elements
>>> np.real_if_close(c,tol=1000) #Return a real array if complex parts close to 0
>>> np.cast['f'](np.pi) #Cast object to a data type
```

#### Other Useful Functions

```
>>> np.angle(b,deg=True) #Return the angle of the complex argument
>>> g = np.linspace(0,np.pi,num=5) #Create an array of evenly spaced values(number of samples)
>>> g [3:] += np.pi
>>> np.unwrap(g) #Unwrap
>>> np.logspace(0,10,3) #Create an array of evenly spaced values (log scale)
>>> np.select([c<4],[c*x]) #Return values from a list of arrays depending on conditions
>>> misc.factorial(a) #Factorial
>>> misc.comb(10,3,exact=True) #Combine N things taken at k time
>>> misc.central_diff_weights(3) #Weights for N-point central derivative
>>> misc.derivative(myfunc,1.0) #Find the n-th derivative of a function at a point
```

### > Linear Algebra

You'll use the linalg and sparse modules.

Note that scipy.linalg contains and expands on numpy.linalg.

```
>>> from scipy import linalg, sparse
```

#### Creating Matrices

```
>>> A = np.matrix(np.random.random((2,2)))
>>> B = np.asmatrix(b)
>>> C = np.mat(np.random.random((10,5)))
>>> D = np.mat([[3,4], [5,6]])
```

#### Basic Matrix Routines

##### Inverse

```
>>> A.I #Inverse
>>> linalg.inv(A) #Inverse
>>> A.T #Transpose matrix
>>> A.H #Conjugate transposition
>>> np.trace(A) #Trace
```

##### Norm

```
>>> linalg.norm(A) #Frobenius norm
>>> linalg.norm(A,1) #L1 norm (max column sum)
>>> linalg.norm(A,np.inf) #L inf norm (max row sum)
```

##### Rank

```
>>> np.linalg.matrix_rank(C) #Matrix rank
```

##### Determinant

```
>>> linalg.det(A) #Determinant
```

##### Solving linear problems

```
>>> linalg.solve(A,b) #Solver for dense matrices
>>> E = np.mat(a).T #Solver for dense matrices
>>> linalg.lstsq(D,E) #Least-squares solution to linear matrix equation
```

##### Generalized inverse

```
>>> linalg.pinv(C) #Compute the pseudo-inverse of a matrix (least-squares solver)
```

```
>>> linalg.pinv2(C) #Compute the pseudo-inverse of a matrix (SVD)
```

#### Creating Sparse Matrices

```
>>> F = np.eye(3, k=1) #Create a 2X2 identity matrix
>>> G = np.mat(np.identity(2)) #Create a 2x2 identity matrix
>>> C[C > 0.5] = 0
>>> H = sparse.csr_matrix(C) #Compressed Sparse Row matrix
>>> I = sparse.csc_matrix(D) #Compressed Sparse Column matrix
>>> J = sparse.dok_matrix(A) #Dictionary Of Keys matrix
>>> E.todense() #Sparse matrix to full matrix
>>> sparse.isspmatrix_csc(A) #Identify sparse matrix
```

#### Sparse Matrix Routines

##### Inverse

```
>>> sparse.linalg.inv(I) #Inverse
```

##### Norm

```
>>> sparse.linalg.norm(I) #Norm
```

##### Solving linear problems

```
>>> sparse.linalg.spsolve(H,I) #Solver for sparse matrices
```

#### Sparse Matrix Functions

```
>>> sparse.linalg.expm(I) #Sparse matrix exponential
```

#### Sparse Matrix Decompositions

```
>>> la, v = sparse.linalg.eigs(F,1) #Eigenvalues and eigenvectors
>>> sparse.linalg.svds(H, 2) #SVD
```

### Matrix Functions

##### Addition

```
>>> np.add(A,D) #Addition
```

##### Subtraction

```
>>> np.subtract(A,D) #Subtraction
```

##### Division

```
>>> np.divide(A,D) #Division
```

##### Multiplication

```
>>> np.multiply(D,A) #Multiplication
>>> np.dot(A,D) #Dot product
>>> np.vdot(A,D) #Vector dot product
>>> np.inner(A,D) #Inner product
>>> np.outer(A,D) #Outer product
>>> np.tensordot(A,D) #Tensor dot product
>>> np.kron(A,D) #Kronecker product
```

##### Exponential Functions

```
>>> linalg.expm(A) #Matrix exponential
>>> linalg.expm2(A) #Matrix exponential (Taylor Series)
>>> linalg.expm3(D) #Matrix exponential (eigenvalue decomposition)
```

##### Logarithm Function

```
>>> linalg.logm(A) #Matrix logarithm
```

##### Trigonometric Functions

```
>>> linalg.sinm(D) #Matrix sine
>>> linalg.cosm(D) #Matrix cosine
>>> linalg.tanm(A) #Matrix tangent
```

##### Hyperbolic Trigonometric Functions

```
>>> linalg.sinhm(D) #Hyperbolic matrix sine
>>> linalg.coshm(D) #Hyperbolic matrix cosine
>>> linalg.tanhm(A) #Hyperbolic matrix tangent
```

##### Matrix Sign Function

```
>>> np.sign(A) #Matrix sign function
```

##### Matrix Square Root

```
>>> linalg.sqrtm(A) #Matrix square root
```

##### Arbitrary Functions

```
>>> linalg.funm(A, lambda x: x*x) #Evaluate matrix function
```

### Decompositions

##### Eigenvalues and Eigenvectors

```
>>> la, v = linalg.eig(A) #Solve ordinary or generalized eigenvalue problem for square matrix
>>> l1, l2 = la #Unpack eigenvalues
>>> v[:,0] #First eigenvector
>>> v[:,1] #Second eigenvector
>>> linalg.eigvals(A) #Unpack eigenvalues
```

##### Singular Value Decomposition

```
>>> U,s,Vh = linalg.svd(B) #Singular Value Decomposition (SVD)
>>> M,N = B.shape
>>> Sig = linalg.diagsvd(s,M,N) #Construct sigma matrix in SVD
```

##### LU Decomposition

```
>>> P,L,U = linalg.lu(C) #LU Decomposition
```

### > Asking For Help

```
>>> help(scipy.linalg.diagsvd)
>>> np.info(np.matrix)
```

Learn Data Skills Online at  
[www.DataCamp.com](http://www.DataCamp.com)

# Python For Data Science

## Importing Data Cheat Sheet

Learn Python online at [www.DataCamp.com](http://www.DataCamp.com)

### > Importing Data in Python

Most of the time, you'll use either NumPy or pandas to import your data:

```
>>> import numpy as np
>>> import pandas as pd
```

### > Help

```
>>> np.info(np.ndarray.dtype)
>>> help(pd.read_csv)
```

### > Text Files

#### Plain Text Files

```
>>> filename = 'huck_finn.txt'
>>> file = open(filename, mode='r') #Open the file for reading
>>> text = file.read() #Read a file's contents
>>> print(file.closed) #Check whether file is closed
>>> file.close() #Close file
>>> print(text)

Using the context manager with

>>> with open('huck_finn.txt', 'r') as file:
    print(file.readline()) #Read a single line
    print(file.readline())
    print(file.readline())
    print(file.readline())
```

#### Table Data: Flat Files

##### Importing Flat Files with NumPy

```
>>> filename = 'huck_finn.txt'
>>> file = open(filename, mode='r') #Open the file for reading
>>> text = file.read() #Read a file's contents
>>> print(file.closed) #Check whether file is closed
>>> file.close() #Close file
>>> print(text)
```

##### Files with one data type

```
>>> filename = 'mnist.txt'
>>> data = np.loadtxt(filename,
    delimiter=',', #String used to separate values
    skiprows=2, #Skip the first 2 lines
    usecols=[0,2], #Read the 1st and 3rd column
    dtype=str) #The type of the resulting array
```

##### Files with mixed data type

```
>>> filename = 'titanic.csv'
>>> data = np.genfromtxt(filename,
    delimiter=',',
    names=True, #Look for column header
    dtype=None)
>>> data_array = np.recfromcsv(filename)
#The default dtype of the np.recfromcsv() function is None
```

##### Importing Flat Files with Pandas

```
>>> filename = 'winequality-red.csv'
>>> data = pd.read_csv(filename,
    nrows=5, #Number of rows of file to read
    header=None, #Row number to use as col names
    sep='\t', #Delimiter to use
    comment='#', #Character to split comments
    na_values=['']) #String to recognize as NA/Nan
```

### > Exploring Your Data

#### NumPy Arrays

```
>>> data_array.dtype #Data type of array elements
>>> data_array.shape #Array dimensions
>>> len(data_array) #Length of array
```

#### Pandas DataFrames

```
>>> df.head() #Return first DataFrame rows
>>> df.tail() #Return last DataFrame rows
>>> df.index #Describe index
>>> df.columns #Describe DataFrame columns
>>> df.info() #Info on DataFrame
>>> data_array = data.values #Convert a DataFrame to an a NumPy array
```

### > SAS File

```
>>> from sas7bdat import SAS7BDAT
>>> with SAS7BDAT('urbanpop.sas7bdat') as file:
    df_sas = file.to_data_frame()
```

### > Stata File

```
>>> data = pd.read_stata('urbanpop.dta')
```

### > Excel Spreadsheets

```
>>> file = 'urbanpop.xlsx'
>>> data = pd.ExcelFile(file)
>>> df_sheet2 = data.parse('1960-1966',
    skiprows=[0],
    names=['Country',
    'AAM: War(2002)'])
>>> df_sheet1 = data.parse(0,
    parse_cols=[0],
    skiprows=[0],
    names=['Country'])

To access the sheet names, use the sheet_names attribute:
>>> data.sheet_names
```

### > Relational Databases

```
>>> from sqlalchemy import create_engine
>>> engine = create_engine('sqlite:///Northwind.sqlite')
Use the table_names() method to fetch a list of table names:
>>> table_names = engine.table_names()
```

#### Querying Relational Databases

```
>>> con = engine.connect()
>>> rs = con.execute("SELECT * FROM Orders")
>>> df = pd.DataFrame(rs.fetchall())
>>> df.columns = rs.keys()
>>> con.close()
```

##### Using the context manager with

```
>>> with engine.connect() as con:
    rs = con.execute("SELECT OrderID FROM Orders")
    df = pd.DataFrame(rs.fetchmany(size=5))
    df.columns = rs.keys()
```

#### Querying relational databases with pandas

```
>>> df = pd.read_sql_query("SELECT * FROM Orders", engine)
```

### > Pickled Files

```
>>> import pickle
>>> with open('pickled_fruit.pkl', 'rb') as file:
    pickled_data = pickle.load(file)
```

### > Matlab Files

```
>>> import scipy.io
>>> filename = 'workspace.mat'
>>> mat = scipy.io.loadmat(filename)
```

### > HDF5 Files

```
>>> import h5py
>>> filename = 'H-H1_LOSC_4_v1-815411200-4096.hdf5'
>>> data = h5py.File(filename, 'r')
```

### > Exploring Dictionaries

Querying relational databases with pandas

```
>>> print(mat.keys()) #Print dictionary keys
>>> for key in data.keys(): #Print dictionary keys
    print(key)
meta
quality
strain
>>> pickled_data.values() #Return dictionary values
>>> print(mat.items()) #Returns items in list format of (key, value) tuple pairs
```

#### Accessing Data Items with Keys

```
>>> for key in data['meta'].keys(): #Explore the HDF5
    structure
        print(key)
Description
DescriptionURL
Detector
Duration
GPSstart
Observatory
Type
UTCstart
#Retrieve the value for a key
>>> print(data['meta']['Description'].value)
```

### > Navigating Your FileSystem

#### Magic Commands

```
!ls #List directory contents of files and directories
%cd .. #Change current working directory
%pwd #Return the current working directory path
```

#### OS Library

```
>>> import os
>>> path = "/usr/tmp"
>>> wd = os.getcwd() #Store the name of current directory in a string
>>> os.listdir(wd) #Output contents of the directory in a list
>>> os.chdir(path) #Change current working directory
>>> os.rename("test1.txt", "test2.txt") #Rename a file
    "test2.txt")
>>> os.remove("test1.txt") #Delete an existing file
>>> os.mkdir("newdir") #Create a new directory
```

# Python For Data Science

## Pandas Basics Cheat Sheet

Learn Pandas Basics online at [www.DataCamp.com](http://www.DataCamp.com)

### Pandas

The **Pandas** library is built on NumPy and provides easy-to-use **data structures** and **data analysis** tools for the Python programming language.

Use the following import convention:

```
>>> import pandas as pd
```

### Pandas Data Structures

#### Series

A **one-dimensional** labeled array capable of holding any data type

a	3
b	-5
c	7
d	4

Index →

```
>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```

#### Dataframe

A **two-dimensional** labeled data structure with columns of potentially different types

	Country	Capital	Population
Index →	Belgium	Brussels	11190846
0	India	New Delhi	1303171035
1	Brazil	Brasilia	207847528

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],
   'Capital': ['Brussels', 'New Delhi', 'Brasilia'],
   'Population': [11190846, 1303171035, 207847528]}
>>> df = pd.DataFrame(data,
   columns=['Country', 'Capital', 'Population'])
```

### Dropping

```
>>> s.drop(['a', 'c']) #Drop values from rows (axis=0)
>>> df.drop('Country', axis=1) #Drop values from columns(axis=1)
```

### Asking For Help

```
>>> help(pd.Series.loc)
```

### Sort & Rank

```
>>> df.sort_index() #Sort by labels along an axis
>>> df.sort_values(by='Country') #Sort by the values along an axis
>>> df.rank() #Assign ranks to entries
```

### > I/O

#### Read and Write to CSV

```
>>> pd.read_csv('file.csv', header=None, nrows=5)
>>> df.to_csv('myDataFrame.csv')
```

#### Read and Write to Excel

```
>>> pd.read_excel('file.xlsx')
>>> df.to_excel('dir/myDataFrame.xlsx', sheet_name='Sheet1')

Read multiple sheets from the same file
>>> xlsx = pd.ExcelFile('file.xlsx')
>>> df = pd.read_excel(xlsx, 'Sheet1')
```

#### Read and Write to SQL Query or Database Table

```
>>> from sqlalchemy import create_engine
>>> engine = create_engine('sqlite:///memory:')
>>> pd.read_sql("SELECT * FROM my_table;", engine)
>>> pd.read_sql_table('my_table', engine)
>>> pd.read_sql_query("SELECT * FROM my_table;", engine)

read_sql() is a convenience wrapper around read_sql_table() and read_sql_query()
>>> df.to_sql('myDF', engine)
```

### > Selection

Also see NumPy Arrays

#### Getting

```
>>> s['b'] #Get one element
-5
>>> df[1:] #Get subset of a DataFrame
   Country Capital Population
1 India New Delhi 1303171035
2 Brazil Brasilia 207847528
```

#### Selecting, Boolean Indexing & Setting

##### By Position

```
>>> df.iloc[[0],[0]] #Select single value by row & column
'Belgium'
>>> df.iat[[0],[0]]
'Belgium'
```

##### By Label

```
>>> df.loc[[0], ['Country']] #Select single value by row & column labels
'Belgium'
>>> df.at[[0], ['Country']]
'Belgium'
```

##### By Label/Position

```
>>> df.ix[2] #Select single row of subset of rows
Country Brazil
Capital Brasilia
Population 207847528
>>> df.ix[:, 'Capital'] #Select a single column of subset of columns
0 Brussels
1 New Delhi
2 Brasilia
>>> df.ix[1, 'Capital'] #Select rows and columns
'New Delhi'

Boolean Indexing
```

```
>>> s[~(s > 1)] #Series s where value is not >1
>>> s[(s < -1) | (s > 2)] #s where value is <-1 or >2
>>> df[df['Population']>1200000000] #Use filter to adjust DataFrame
```

##### Setting

```
>>> s['a'] = 6 #Set index a of Series s to 6
```

### > Retrieving Series/DataFrame Information

#### Basic Information

```
>>> df.shape #(rows,columns)
>>> df.index #Describe index
>>> df.columns #Describe DataFrame columns
>>> df.info() #Info on DataFrame
>>> df.count() #Number of non-NA values
```

#### Summary

```
>>> df.sum() #Sum of values
>>> df.cumsum() #Cumulative sum of values
>>> df.min()/df.max() #Minimum/maximum values
>>> df.idxmin()/df.idxmax() #Minimum/Maximum index value
>>> df.describe() #Summary statistics
>>> df.mean() #Mean of values
>>> df.median() #Median of values
```

### > Applying Functions

```
>>> f = lambda x: x*2
>>> df.apply(f) #Apply function
>>> df.applymap(f) #Apply function element-wise
```

### > Data Alignment

#### Internal Data Alignment

NA values are introduced in the indices that don't overlap:

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
>>> s + s3
a 10.0
b NaN
c 5.0
d 7.0
```

#### Arithmetic Operations with Fill Methods

You can also do the internal data alignment yourself with the help of the fill methods:

```
>>> s.add(s3, fill_value=0)
a 10.0
b -5.0
c 5.0
d 7.0
>>> s.sub(s3, fill_value=2)
>>> s.div(s3, fill_value=4)
>>> s.mul(s3, fill_value=3)
```

Learn Data Skills Online at  
[www.DataCamp.com](http://www.DataCamp.com)

# Python For Data Science

## Data Wrangling in Pandas Cheat Sheet

Learn Data Wrangling online at [www.DataCamp.com](http://www.DataCamp.com)

## > Reshaping Data

### Pivot

```
>>> df3 = df2.pivot(index='Date', #Spread rows into columns
                   columns='Type',
                   values='Value')
```

	Date	Type	Value
0	2016-03-01	a	11.432
1	2016-03-02	b	13.031
2	2016-03-01	c	20.784
3	2016-03-03	a	99.906
4	2016-03-02	a	1.303
5	2016-03-03	c	20.784

	Type	a	b	c
2016-03-01		11.432		
2016-03-02			13.031	
2016-03-03		99.906		20.784

### Pivot Table

```
>>> df4 = pd.pivot_table(df2, #Spread rows into
                        columns values='Value',
                        index='Date',
                        columns='Type'))
```

### Stack / Unstack

```
>>> stacked = df5.stack() #Pivot a level of column labels
>>> stacked.unstack() #Pivot a level of index labels
```

	0	1
1	0.233482	0.390959
2	0.184713	0.237102
3	0.433522	0.429401

	0	1
1	0.233482	0.390959
2	0.184713	0.237102
3	0.433522	0.429401

Unstacked

Stacked

### Melt

```
>>> pd.melt(df2, #Gather columns into rows
            id_vars=['Date'],
            value_vars=["Type", "Value"],
            value_name="Observations")
```

	Date	Type	Value
0	2016-03-01	a	11.432
1	2016-03-02	b	13.031
2	2016-03-01	c	20.784
3	2016-03-03	a	99.906
4	2016-03-02	a	1.303
5	2016-03-03	c	20.784

	Date	Type	Value
0	2016-03-01	a	11.432
1	2016-03-02	b	13.031
2	2016-03-01	c	20.784
3	2016-03-03	a	99.906
4	2016-03-02	a	1.303
5	2016-03-03	c	20.784

## > Iteration

```
>>> df.iteritems() #(Column-index, Series) pairs
>>> df.iterrows() #(Row-index, Series) pairs
```

## > Missing Data

```
>>> df.dropna() #Drop NaN values
>>> df3.fillna(df3.mean()) #Fill NaN values with a predetermined value
>>> df2.replace("a", "f") #Replace values with others
```

## > Advanced Indexing

Also see NumPy Arrays

### Selecting

```
>>> df3.loc[:, (df3>1).any()] #Select cols with any vals >1
>>> df3.loc[:, (df3>1).all()] #Select cols with vals > 1
>>> df3.loc[:, df3.isnull().any()] #Select cols with NaN
>>> df3.loc[:, df3.notnull().all()] #Select cols without NaN
```

### Indexing With isin()

```
>>> df[(df.Country.isin(df2.Type))] #Find same elements
>>> df3.filter(items="a", "b") #Filter on values
>>> df.select(lambda x: not x%5) #Select specific elements
```

### Where

```
>>> s.where(s > 0) #Subset the data
```

### Query

```
>>> df6.query('second > first') #Query DataFrame
```

## Setting/Resetting Index

```
>>> df.set_index('Country') #Set the index
>>> df4 = df.reset_index() #Reset the index
>>> df = df.rename(index=str, #Rename
                  DataFrame columns={"Country":"cntry",
                                     "Capital":"cptl",
                                     "Population":"ppltn"})
```

## Reindexing

```
>>> s2 = s.reindex(['a', 'c', 'd', 'e', 'b'])
```

### Forward Filling

```
>>> df.reindex(range(4),
               method='ffill')
```

Country	Capital	Population
0 Belgium	Brussels	11190846
1 India	New Delhi	1303171035
2 Brazil	Brasilia	207847528
3 Brazil	Brasilia	207847528

### Backward Filling

```
>>> s3 = s.reindex(range(5),
                   method='bfill')
```

	0	3
1	3	3
2	3	3
3	3	3
4	3	3

## MultilIndexing

```
>>> arrays = [np.array([1,2,3]),
             np.array([5,4,3])]
>>> df5 = pd.DataFrame(np.random.rand(3, 2), index=arrays)
>>> tuples = list(zip(*arrays))
>>> index = pd.MultiIndex.from_tuples(tuples,
                                       names=['first', 'second'])
>>> df6 = pd.DataFrame(np.random.rand(3, 2), index=index)
>>> df2.set_index(["Date", "Type"])
```

## > Duplicate Data

```
>>> s3.unique() #Return unique values
>>> df2.duplicated('Type') #Check duplicates
>>> df2.drop_duplicates('Type', keep='last') #Drop duplicates
>>> df.index.duplicated() #Check index duplicates
```

## > Grouping Data

### Aggregation

```
>>> df2.groupby(by=['Date', 'Type']).mean()
>>> df4.groupby(level=0).sum()
>>> df4.groupby(level=0).agg({'a':lambda x:sum(x)/len(x), 'b': np.sum})
```

### Transformation

```
>>> customSum = lambda x: (x+x%2)
>>> df4.groupby(level=0).transform(customSum)
```

## > Combining Data

data1		data2	
X1	X2	X1	X3
a			a
b			b
c			NaN

### Merge

```
>>> pd.merge(data1,
             data2,
             how='left',
             on='X1')
```

X1		X2	X3
a		11.432	20.784
b		1.303	NaN
c		99.906	NaN

```
>>> pd.merge(data1,
             data2,
             how='right',
             on='X1')
```

X1		X2	X3
a		11.432	20.784
b		1.303	NaN
c		NaN	20.784

```
>>> pd.merge(data1,
             data2,
             how='inner',
             on='X1')
```

X1		X2	X3
a		11.432	20.784
b		1.303	NaN

```
>>> pd.merge(data1,
             data2,
             how='outer',
             on='X1')
```

X1		X2	X3
a		11.432	20.784
b		1.303	NaN
c		99.906	NaN
d		NaN	20.784

# The Data Visualization Cheat Sheet

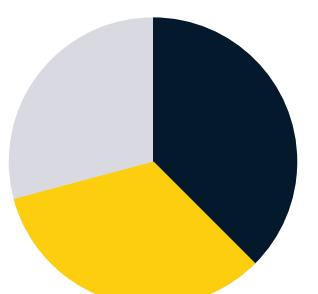
Learn Data Visualization online at [www.DataCamp.com](http://www.DataCamp.com)

## How to use this cheat sheet

Use this cheat sheet for inspiration when making your next data visualizations. For more data visualization cheat sheets, check out our cheat sheets repository [here](#).

## > Part-to-whole charts

Pie chart



One of the most common ways to show part to whole data. It is also commonly used with percentages

Donut pie chart



The donut pie chart is a variant of the pie chart, the difference being it has a hole in the center for readability

Heat maps



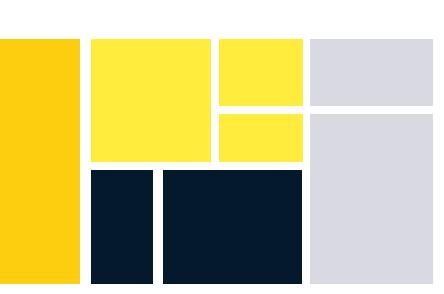
Heatmaps are two-dimensional charts that use color shading to represent data trends.

Stacked column chart



Best to compare subcategories within categorical data. Can also be used to compare percentages

Treemap charts



2D rectangles whose size is proportional to the value being measured and can be used to display hierarchically structured data

USE CASES

- 1. Voting preference by age group
- 2. Market share of cloud providers

USE CASES

- 1. Android OS market share
- 2. Monthly sales by channel

USE CASES

- 1. Average monthly temperatures across the year
- 2. Departments with the highest amount of attrition over time

USE CASES

- 1. Quarterly sales per region
- 2. Total car sales by producer

USE CASES

- 1. Grocery sales count with categories
- 2. Stock price comparison by industry and company

## > Capture a trend

Line chart



The most straightforward way to capture how a numeric variable is changing over time

Multi-line chart



Captures multiple numeric variables over time. It can include multiple axes allowing comparison of different units and scale ranges

Area chart



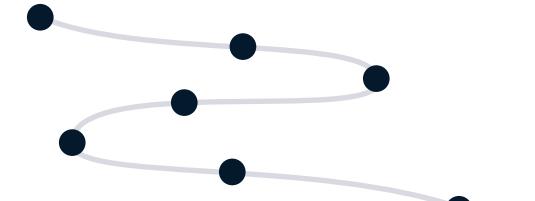
Shows how a numeric value progresses by shading the area between line and the x-axis

Stacked area chart



Most commonly used variation of area charts, the best use is to track the breakdown of a numeric value by subgroups

Spline chart



Smoothed version of a line chart. It differs in that data points are connected with smoothed curves to account for missing values, as opposed to straight lines

USE CASES

- 1. Revenue in \$ over time
- 2. Energy consumption in kWh over time
- 3. Google searches over time

USE CASES

- 1. Apple vs Amazon stocks over time
- 2. LeBron vs Steph Curry searches over time
- 3. Bitcoin vs Ethereum price over time

USE CASES

- 1. Total sales over time
- 2. Active users over time

USE CASES

- 1. Active users over time by segment
- 2. Total revenue over time by country

USE CASES

- 1. Electricity consumption over time
- 2. CO2 emissions over time

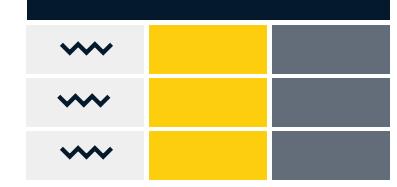
## > Visualize a single value

Card



Cards are great for showing and tracking KPIs in dashboards or presentations

Table chart



Best to be used on small datasets, it displays tabular data in a table

Gauge chart



This chart is often used in executive dashboard reports to show relevant KPIs

USE CASES

- 1. Revenue to date on a sales dashboard
- 2. Total sign-ups after a promotion

USE CASES

- 1. Account executive leaderboards
- 2. Registrations per webinar

USE CASES

- 1. NPS scores
- 2. Revenue to target

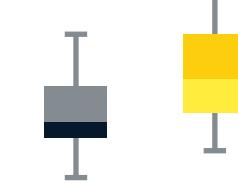
## > Capture distributions

Histogram



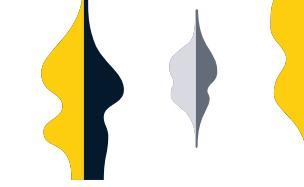
Shows the distribution of a variable. It converts numerical data into bins as columns. The x-axis shows the range, and the y-axis represents the frequency

Box plot



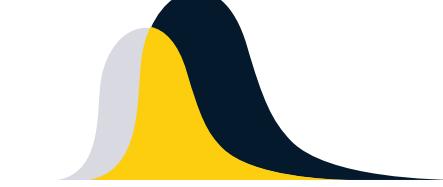
Shows the distribution of a variable using 5 key summary statistics—minimum, first quartile, median, third quartile, and maximum

Violin plot



A variation of the box plot. It also shows the full distribution of the data alongside summary statistics

Density plot



Visualizes a distribution by using smoothing to allow smoother distributions and better capture the distribution shape of the data

## > Visualize relationships

Bar chart



One of the easiest charts to read which helps in quick comparison of categorical data. One axis contains categories and the other axis represents values

Column chart



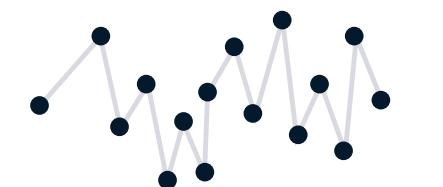
Also known as a vertical bar chart, where the categories are placed on the x-axis. These are preferred over bar charts for short labels, date ranges, or negatives in values

Scatter plot



Most commonly used chart when observing the relationship between two variables. It is especially useful for quickly surfacing potential correlations between data points

Connected scatterplot



A hybrid between a scatter plot and a line plot, the scatter dots are connected with a line

Bubble chart



Often used to visualize data points with 3 dimensions, namely visualized on the x-axis, y-axis, and with the size of the bubble. It tries to show relations between data points using location and size

Word cloud chart



A convenient visualization for visualizing the most prevalent words that appear in a text

## > Visualize a flow

Sankey chart



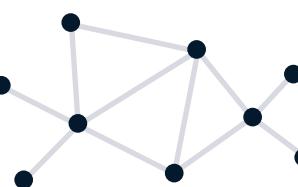
Useful for representing flows in systems. This flow can be any measurable quantity

Chord chart



Useful for presenting weighted relationships or flows between nodes. Especially useful for highlighting the dominant or important flows

Network chart



Similar to a graph, it consists of nodes and interconnected edges. It illustrates how different items have relationships with each other

Learn Data Skills Online at  
[www.DataCamp.com](http://www.DataCamp.com)

# Python For Data Science

## Matplotlib Cheat Sheet

Learn Matplotlib online at [www.DataCamp.com](http://www.DataCamp.com)

## Matplotlib

Matplotlib is a Python 2D plotting library which produces publication-quality figures in a variety of hardcopy formats and interactive environments across platforms.

## Prepare The Data

### 1D Data

```
>>> import numpy as np
>>> x = np.linspace(0, 10, 100)
>>> y = np.cos(x)
>>> z = np.sin(x)
```

### 2D Data or Images

```
>>> data = 2 * np.random.random((10, 10))
>>> data2 = 3 * np.random.random((10, 10))
>>> Y, X = np.mgrid[-3:3:100j, -3:3:100j]
>>> U = -1 - X**2 + Y
>>> V = 1 + X - Y**2
>>> from matplotlib.cbook import get_sample_data
>>> img = np.load(get_sample_data('axes_grid/bivariate_normal.npy'))
```

## Create Plot

```
>>> import matplotlib.pyplot as plt
```

### Figure

```
>>> fig = plt.figure()
>>> fig2 = plt.figure(figsize=plt.figaspect(2.0))
```

### Axes

All plotting is done with respect to an Axes. In most cases, a subplot will fit your needs. A subplot is an axes on a grid system.

```
>>> fig.add_axes()
>>> ax1 = fig.add_subplot(221) #row-col-num
>>> ax3 = fig.add_subplot(212)
>>> fig3, axes = plt.subplots(nrows=2, ncols=2)
>>> fig4, axes2 = plt.subplots(ncols=3)
```

## Save Plot

```
>>> plt.savefig('foo.png') #Save figures
>>> plt.savefig('foo.png', transparent=True) #Save transparent figures
```

## Show Plot

```
>>> plt.show()
```

## Plotting Routines

### 1D Data

```
>>> fig, ax = plt.subplots()
>>> lines = ax.plot(x,y) #Draw points with lines or markers connecting them
>>> ax.scatter(x,y) #Draw unconnected points, scaled or colored
>>> axes[0,0].bar([1,2,3],[3,4,5]) #Plot vertical rectangles (constant width)
>>> axes[0,0].barh([0.5,1,2.5],[0,1,2]) #Plot horizontal rectangles (constant height)
>>> axes[1,1].axhline(0.45) #Draw a horizontal line across axes
>>> axes[0,1].axvline(0.65) #Draw a vertical line across axes
>>> ax.fill(x,y,color='blue') #Draw filled polygons
>>> ax.fill_between(x,y,color='yellow') #Fill between y-values and 0
```

### 2D Data

```
>>> fig, ax = plt.subplots()
>>> im = ax.imshow(img, #Colormapped or RGB arrays
                  cmap='gist_earth',
                  interpolation='nearest',
                  vmin=-2,
                  vmax=2)
>>> axes2[0].pcolor(data2) #Pseudocolor plot of 2D array
>>> axes2[0].pcolormesh(data) #Pseudocolor plot of 2D array
>>> CS = plt.contour(Y,X,U) #Plot contours
>>> axes2[2].contourf(data1) #Plot filled contours
>>> axes2[2]= ax.clabel(CS) #Label a contour plot
```

### Vector Fields

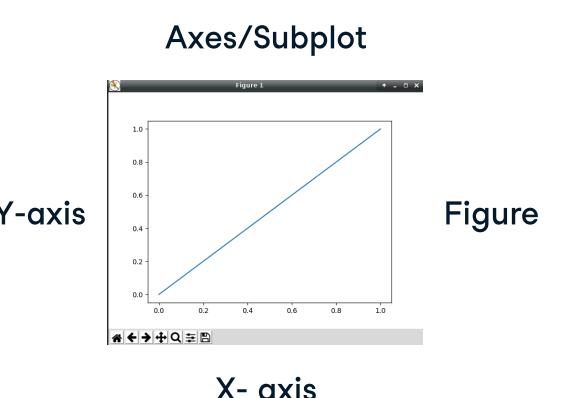
```
>>> axes[0,1].arrow(0,0,0.5,0.5) #Add an arrow to the axes
>>> axes[1,1].quiver(y,z) #Plot a 2D field of arrows
>>> axes[0,1].streamplot(X,Y,U,V) #Plot a 2D field of arrows
```

### Data Distributions

```
>>> ax1.hist(y) #Plot a histogram
>>> ax3.boxplot(y) #Make a box and whisker plot
>>> ax3.violinplot(z) #Make a violin plot
```

## Plot Anatomy & Workflow

### Plot Anatomy



### Workflow

The basic steps to creating plots with matplotlib are:

- 1 Prepare Data
  - 2 Create Plot
  - 3 Plot
  - 4 Customized Plot
  - 5 Save Plot
  - 6 Show Plot
- ```
>>> import matplotlib.pyplot as plt
>>> x = [1,2,3,4] #Step 1
>>> y = [10,20,25,30]
>>> fig = plt.figure() #Step 2
>>> ax = fig.add_subplot(111) #Step 3
>>> ax.plot(x, y, color='lightblue', linewidth=3) #Step 3, 4
>>> ax.scatter([2,4,6],
              [5,15,25],
              color='darkgreen',
              marker='^')
>>> ax.set_xlim(1, 6.5)
>>> plt.savefig('foo.png') #Step 5
>>> plt.show() #Step 6
```

## Close and Clear

```
>>> plt.cla() #Clear an axis
>>> plt.clf() #Clear the entire figure
>>> plt.close() #Close a window
```

## Plotting Cutomize Plot

### Colors, Color Bars & Color Maps

```
>>> plt.plot(x, x, x*x2, x, x*x3)
>>> ax.plot(x, y, alpha = 0.4)
>>> ax.plot(x, y, c='k')
>>> fig.colorbar(im, orientation='horizontal')
>>> im = ax.imshow(img,
                  cmap='seismic')
```

### Markers

```
>>> fig, ax = plt.subplots()
>>> ax.scatter(x,y,marker=".")
>>> ax.plot(x,y,marker="o")
```

### Linestyles

```
>>> plt.plot(x,y,linewidth=4.0)
>>> plt.plot(x,y,ls='solid')
>>> plt.plot(x,y,ls='--')
>>> plt.plot(x,y,'--',x*x2,y**2,'-.')
>>> plt.setp(lines,color='r',linewidth=4.0)
```

### Text & Annotations

```
>>> ax.text(1,
           -2.1,
           'Example Graph',
           style='italic')
>>> ax.annotate("Sine",
               xy=(8, 0),
               xycoords='data',
               xytext=(10.5, 0),
               textcoords='data',
               arrowprops=dict(arrowstyle="→",
                               connectionstyle="arc3"))

```

### MathText

```
>>> plt.title(r'$\sigma_i=15$', fontsize=20)
```

### Limits, Legends and Layouts

#### Limits & Autoscaling

```
>>> ax.margins(x=0.0,y=0.1) #Add padding to a plot
>>> ax.axis('equal') #Set the aspect ratio of the plot to 1
>>> ax.set(xlim=[0,10.5],ylim=[-1.5,1.5]) #Set limits for x-and y-axis
>>> ax.set_xlim(0,10.5) #Set limits for x-axis
```

#### Legends

```
>>> ax.set(title='An Example Axes', #Set a title and x-and y-axis labels
           ylabel='Y-Axis',
           xlabel='X-Axis')
>>> ax.legend(loc='best') #No overlapping plot elements
```

#### Ticks

```
>>> ax.xaxis.set(ticks=range(1,5), #Manually set x-ticks
                  ticklabels=[3,100,-12,"foo"])
>>> ax.tick_params(axis='y', #Make y-ticks longer and go in and out
                  direction='inout',
                  length=10)
```

#### Subplot Spacing

```
>>> fig3.subplots_adjust(wspace=0.5, #Adjust the spacing between subplots
                           hspace=0.3,
                           left=0.125,
                           right=0.9,
                           top=0.9,
                           bottom=0.1)
>>> fig.tight_layout() #Fit subplot(s) in to the figure area
```

#### Axis Spines

```
>>> ax1.spines['top'].set_visible(False) #Make the top axis line for a plot invisible
>>> ax1.spines['bottom'].set_position(('outward',10)) #Move the bottom axis line outward
```

# Python For Data Science Seaborn Cheat Sheet

Learn Seaborn online at [www.DataCamp.com](http://www.DataCamp.com)

## Statistical Data Visualization With Seaborn

The Python visualization library **Seaborn** is based on `matplotlib` and provides a high-level interface for drawing attractive statistical graphics.

**Make use of the following aliases to import the libraries:**

```
>>> import matplotlib.pyplot as plt
>>> import seaborn as sns
```

The basic steps to creating plots with Seaborn are:

1. Prepare some data
2. Control figure aesthetics
3. Plot with Seaborn
4. Further customize your plot
5. Show your plot

```
>>> import matplotlib.pyplot as plt
>>> import seaborn as sns
>>> tips = sns.load_dataset("tips") #Step 1
>>> sns.set_style("whitegrid") #Step 2
>>> g = sns.lmplot(x="tip", #Step 3
                   y="total_bill",
                   data=tips,
                   aspect=2)
>>> g.set_axis_labels("Tip","Total bill(USD)").set(xlim=(0,10),ylim=(0,100))
>>> plt.title("title") #Step 4
>>> plt.show(g) #Step 5
```

## 1 Data

Also see [Lists, NumPy & Pandas](#)

```
>>> import pandas as pd
>>> import numpy as np
>>> uniform_data = np.random.rand(10, 12)
>>> data = pd.DataFrame({‘x’:np.arange(1,101),
                       ‘y’:np.random.normal(0,4,100)})
```

Seaborn also offers built-in data sets:

```
>>> titanic = sns.load_dataset("titanic")
>>> iris = sns.load_dataset("iris")
```

## 2 Figure Aesthetics

Also see [Matplotlib](#)

```
>>> f, ax = plt.subplots(figsize=(5,6)) #Create a figure and one subplot
```

### Seaborn styles

```
>>> sns.set() #Re)set the seaborn default
>>> sns.set_style("whitegrid") #Set the matplotlib parameters
>>> sns.set_style("ticks", #Set the matplotlib parameters
                  {“x tick.major.size”:8,
                   “y tick.major.size”:8})
#Return a dict of params or use with to temporarily set the style
>>> sns.axes_style("whitegrid")
```

## 3 Plotting With Seaborn

### Axis Grids

```
>>> g = sns.FacetGrid(titanic, #Subplot grid for plotting conditional relationships
                      col="survived",
                      row="sex")
>>> g.map(plt.hist, "age")
>>> sns.factorplot(x="pclass", #Draw a categorical plot onto a Facetgrid
                     y="survived",
                     hue="sex",
                     data=titanic)
>>> sns.lmplot(x="sepal_width", #Plot data and regression model fits across a FacetGrid
                  y="sepal_length",
                  hue="species",
                  data=iris)
>>> h = sns.PairGrid(iris) #Subplot grid for plotting pairwise relationships
>>> h = h.map(plt.scatter)
>>> sns.pairplot(iris) #Plot pairwise bivariate distributions
>>> i = sns.JointGrid(x="x", #Grid for bivariate plot with marginal univariate plots
                      y="y",
                      data=data)
>>> i = i.plot(sns.regplot,
                  sns.distplot)
>>> sns.jointplot("sepal_length", #Plot bivariate distribution
                  "sepal_width",
                  data=iris,
                  kind='kde')
```

## 4 Further Customizations

Also see [Matplotlib](#)

### Axisgrid Objects

```
>>> g.despine(left=True) #Remove left spine
>>> g.set_ylabels("Survived") #Set the labels of the y-axis
>>> g.set_xticklabels(rotation=45) #Set the tick labels for x
>>> g.set_axis_labels("Survived", #Set the axis labels
                      "Sex")
>>> h.set(xlim=(0,5), #Set the limit and ticks of the x-and y-axis
          ylim=(0,5),
          xticks=[0,2.5,5],
          yticks=[0,2.5,5])
```

### Plot

```
>>> plt.title("A Title") #Add plot title
>>> plt.ylabel("Survived") #Adjust the label of the y-axis
>>> plt.xlabel("Sex") #Adjust the label of the x-axis
>>> plt.ylim(0,100) #Adjust the limits of the y-axis
>>> plt.xlim(0,10) #Adjust the limits of the x-axis
>>> plt.setp(ax,yticks=[0,5]) #Adjust a plot property
>>> plt.tight_layout() #Adjust subplot params
```

### Context Functions

```
>>> sns.set_context("talk") #Set context to "talk"
>>> sns.set_context("notebook", #Set context to "notebook",
                     font_scale=1.5, #Scale font elements and
                     rc={“lines.linewidth”:2.5}) #override param mapping
```

### Color Palette

```
>>> sns.set_palette("husl",3) #Define the color palette
>>> sns.color_palette("husl") #Use with with to temporarily set palette
>>> flatui = [“#9b59b6”, “#3498db”, “#95a5a6”, “#e74c3c”, “#34495e”, “#2ecc71”]
>>> sns.set_palette(flatui) #Set your own color palette
```

### Regression Plots

```
>>> sns.regplot(x="sepal_width", #Plot data and a linear regression model fit
                  y="sepal_length",
                  data=iris,
                  ax=ax)
```

### Distribution Plots

```
>>> plot = sns.distplot(data.y, #Plot univariate distribution
                           kde=False,
                           color="b")
```

### Matrix Plots

```
>>> sns.heatmap(uniform_data,vmin=0,vmax=1) #Heatmap
```

### Categorical Plots

#### Scatterplot

```
>>> sns.stripplot(x="species", #Scatterplot with one categorical variable
                      y="petal_length",
                      data=iris)
>>> sns.swarmplot(x="species", #Categorical scatterplot with non-overlapping points
                      y="petal_length",
                      data=iris)
```

#### Bar Chart

```
>>> sns.barplot(x="sex", #Show point estimates & confidence intervals with scatterplot glyphs
                  y="survived",
                  hue="class",
                  data=titanic)
```

#### Count Plot

```
>>> sns.countplot(x="deck", #Show count of observations
                      data=titanic,
                      palette="Greens_d")
```

#### Point Plot

```
>>> sns.pointplot(x="class", #Show point estimates & confidence intervals as rectangular bars
                      y="survived",
                      hue="sex",
                      data=titanic,
                      palette={"male":“g”,
                               “female”:“m”},
                      markers=[“^”, “o”],
                      linestyles=[“-”, “--”])
```

#### Boxplot

```
>>> sns.boxplot(x="alive", #Boxplot
                  y="age",
                  hue="adult_male",
                  data=titanic)
>>> sns.boxplot(data=iris,orient="h") #Boxplot with wide-form data
```

#### Violinplot

```
>>> sns.violinplot(x="age", #Violin plot
                      y="sex",
                      hue="survived",
                      data=titanic)
```

## 5 Show or Save Plot

Also see [Matplotlib](#)

```
>>> plt.show() #Show the plot
>>> plt.savefig("foo.png") #Save the plot as a figure
>>> plt.savefig("foo.png", #Save transparent figure
                  transparent=True)
```

## > Close & Clear

Also see [Matplotlib](#)

```
>>> plt.cla() #Clear an axis
>>> plt.clf() #Clear an entire figure
>>> plt.close() #Close a window
```

# Python For Data Science

## PySpark RDD Cheat Sheet

Learn PySpark RDD online at [www.DataCamp.com](http://www.DataCamp.com)

### Spark



PySpark is the Spark Python API that exposes the Spark programming model to Python.

### > Initializing Spark

#### SparkContext

```
>>> from pyspark import SparkContext
>>> sc = SparkContext(master = 'local[2]')
```

#### Inspect SparkContext

```
>>> sc.version #Retrieve SparkContext version
>>> sc.pythonVer #Retrieve Python version
>>> sc.master #Master URL to connect to
>>> str(sc.sparkHome) #Path where Spark is installed on worker nodes
>>> str(sc.sparkUser()) #Retrieve name of the Spark User running SparkContext
>>> sc.appName #Return application name
>>> sc.applicationId #Retrieve application ID
>>> sc.defaultParallelism #Return default level of parallelism
>>> sc.defaultMinPartitions #Default minimum number of partitions for RDDs
```

#### Configuration

```
>>> from pyspark import SparkConf, SparkContext
>>> conf = (SparkConf()
...     .setMaster("local")
...     .setAppName("My app")
...     .set("spark.executor.memory", "1g"))
>>> sc = SparkContext(conf = conf)
```

#### Using The Shell

In the PySpark shell, a special interpreter-aware SparkContext is already created in the variable called `sc`.

```
$ ./bin/spark-shell --master local[2]
$ ./bin/pyspark --master local[4] --py-files code.py
```

Set which master the context connects to with the `--master` argument, and add Python .zip, .egg or .py files to the runtime path by passing a comma-separated list to `--py-files`.

### > Loading Data

#### Parallelized Collections

```
>>> rdd = sc.parallelize([('a',7),('a',2),('b',2)])
>>> rdd2 = sc.parallelize([('a',2),('d',1),('b',1)])
>>> rdd3 = sc.parallelize(range(100))
>>> rdd4 = sc.parallelize([('a',[x,y,z]), ('b',[p,r])])
```

#### External Data

Read either one text file from HDFS, a local file system or any Hadoop-supported file system URI with `textFile()`, or read in a directory of text files with `wholeTextFiles()`

```
>>> textFile = sc.textFile("/my/directory/*.txt")
>>> textFile2 = sc.wholeTextFiles("/my/directory/")
```

### > Retrieving RDD Information

#### Basic Information

```
>>> rdd.getNumPartitions() #List the number of partitions
>>> rdd.count() #Count RDD instances 3
>>> rdd.countByKey() #Count RDD instances by key
defaultdict(<type 'int'>,{'a':2,'b':1})
>>> rdd.countByValue() #Count RDD instances by value
defaultdict(<type 'int'>,{'b':2,('a',2):1,('a',7):1})
>>> rdd.collectAsMap() #Return (key,value) pairs as a dictionary
{'a': 2, 'b': 2}
>>> rdd3.sum() #Sum of RDD elements 4950
>>> sc.parallelize([]).isEmpty() #Check whether RDD is empty
True
```

#### Summary

```
>>> rdd3.max() #Maximum value of RDD elements
99
>>> rdd3.min() #Minimum value of RDD elements
0
>>> rdd3.mean() #Mean value of RDD elements
49.5
>>> rdd3.stdev() #Standard deviation of RDD elements
28.86607004772218
>>> rdd3.variance() #Compute variance of RDD elements
833.25
>>> rdd3.histogram(3) #Compute histogram by bins
([0,33,66,99],[33,33,34])
>>> rdd3.stats() #Summary statistics (count, mean, stdev, max & min)
```

### > Applying Functions

```
#Apply a function to each RDD element
>>> rdd.map(lambda x: x+(x[1],x[0])).collect()
[('a',7,7,'a'),('a',2,2,'a'),('b',2,2,'b')]
#Apply a function to each RDD element and flatten the result
>>> rdd5 = rdd.flatMap(lambda x: x+(x[1],x[0]))
>>> rdd5.collect()
['a',7,7,'a','a',2,2,'a','b',2,2,'b']
#Apply a flatMap function to each (key,value) pair of rdd4 without changing the keys
>>> rdd4.flatMapValues(lambda x: x).collect()
[('a','x'),('a','y'),('a','z'),('b','p'),('b','r')]
```

### > Selecting Data

**Getting**

```
>>> rdd.collect() #Return a list with all RDD elements
[('a', 7), ('a', 2), ('b', 2)]
>>> rdd.take(2) #Take first 2 RDD elements
[('a', 7), ('a', 2)]
>>> rdd.first() #Take first RDD element
('a', 7)
>>> rdd.top(2) #Take top 2 RDD elements
[('b', 2), ('a', 7)]
```

**Sampling**

```
>>> rdd3.sample(False, 0.15, 81).collect() #Return sampled subset of rdd3
[3,4,27,31,40,41,42,43,60,76,79,80,86,97]
```

**Filtering**

```
>>> rdd.filter(lambda x: "a" in x).collect() #Filter the RDD
[('a',7),('a',2)]
>>> rdd5.distinct().collect() #Return distinct RDD values
[('a',2),('b',7)]
>>> rdd.keys().collect() #Return (key,value) RDD's keys
[('a', 'a', 'b')]
```

### > Iterating

```
>>> def g(x): print(x)
>>> rdd.foreach(g) #Apply a function to all RDD elements
('a', 7)
('b', 2)
('a', 2)
```

### > Reshaping Data

#### Reducing

```
>>> rdd.reduceByKey(lambda x,y : x+y).collect() #Merge the rdd values for each key
[('a',9),('b',2)]
>>> rdd.reduce(lambda a, b: a + b) #Merge the rdd values
('a',7,'a',2,'b',2)
```

#### Grouping by

```
>>> rdd3.groupBy(lambda x: x % 2) #Return RDD of grouped values
.mapValues(list)
.collect()
>>> rdd.groupByKey() #Group rdd by key
.mapValues(list)
.collect()
[('a',[7,2]),('b',[2])]
```

#### Aggregating

```
>>> seqOp = (lambda x,y: (x[0]+y,x[1]+1))
>>> combOp = (lambda x,y:(x[0]+y[0],x[1]+y[1]))
#Aggregate RDD elements of each partition and then the results
>>> rdd3.aggregate((0,0),seqOp,combOp)
(4950,100)
#Aggregate values of each RDD key
>>> rdd.aggregateByKey((0,0),seqOp,combOp).collect()
[('a',(9,2)), ('b',(2,1))]
#Aggregate the elements of each partition, and then the results
>>> rdd3.fold(0,add)
4950
#Merge the values for each key
>>> rdd.foldByKey(0, add).collect()
[('a',9),('b',2)]
#Create tuples of RDD elements by applying a function
>>> rdd3.keyBy(lambda x: x*x).collect()
```

### > Mathematical Operations

```
>>> rdd.subtract(rdd2).collect() #Return each rdd value not contained in rdd2
[('b',2),('a',7)]
#Return each (key,value) pair of rdd2 with no matching key in rdd
>>> rdd2.subtractByKey(rdd).collect()
[('d', 1)]
>>> rdd.cartesian(rdd2).collect() #Return the Cartesian product of rdd and rdd2
```

### > Sort

```
>>> rdd2.sortBy(lambda x: x[1]).collect() #Sort RDD by given function
[('d',1),('b',1),('a',2)]
>>> rdd2.sortByKey().collect() #Sort (key, value) RDD by key
[('a',2),('b',1),('d',1)]
```

### > Repartitioning

```
>>> rdd.repartition(4) #New RDD with 4 partitions
>>> rdd.coalesce(1) #Decrease the number of partitions in the RDD to 1
```

### > Saving

```
>>> rdd.saveAsTextFile("rdd.txt")
>>> rdd.saveAsHadoopFile("hdfs://namenodehost/parent/child",
...     'org.apache.hadoop.mapred.TextOutputFormat')
```

### > Stopping SparkContext

```
>>> sc.stop()
```

### > Execution

```
$ ./bin/spark-submit examples/src/main/python/pi.py
```

# Python For Data Science

## PySpark SQL Basics Cheat Sheet

Learn PySpark SQL online at [www.DataCamp.com](http://www.DataCamp.com)

## PySpark & Spark SQL



Spark SQL is Apache Spark's module for working with structured data.

### > Initializing SparkSession

A SparkSession can be used create DataFrame, register DataFrame as tables, execute SQL over tables, cache tables, and read parquet files.

```
>>> from pyspark.sql import SparkSession
>>> spark = SparkSession \
    .builder \
    .appName("Python Spark SQL basic example") \
    .config("spark.some.config.option", "some-value") \
    .getOrCreate()
```

### > Creating DataFrames

#### From RDDs

```
>>> from pyspark.sql.types import *
```

#### Infer Schema

```
>>> sc = spark.sparkContext
>>> lines = sc.textFile("people.txt")
>>> parts = lines.map(lambda l: l.split(","))
>>> people = parts.map(lambda p: Row(name=p[0],age=int(p[1])))
>>> peopledf = spark.createDataFrame(people)
```

#### Specify Schema

```
>>> people = parts.map(lambda p: Row(name=p[0],
    age=int(p[1].strip())))
>>> schemaString = "name age"
>>> fields = [StructField(field_name, StringType(), True) for
field_name in schemaString.split()]
>>> schema = StructType(fields)
>>> spark.createDataFrame(people, schema).show()
```

|  | name     | age |
|--|----------|-----|
|  | Mine     | 28  |
|  | Filip    | 29  |
|  | Jonathan | 30  |

#### From Spark Data Sources

##### JSON

```
>>> df = spark.read.json("customer.json")
>>> df.show()
+-----+-----+
| address|age|firstName|lastName|phoneNumber|
+-----+-----+
|[New York,10021,N...| 25| John | Smith |[212 555-1234,ho...|
|[New York,10021,N...| 25| Jane | Doe |[323 888-1234,ho...|
+-----+-----+
```

```
>>> df2 = spark.read.load("people.json", format="json")
```

##### Parquet files

```
>>> df3 = spark.read.load("users.parquet")
```

##### TXT files

```
>>> df4 = spark.read.text("people.txt")
```

### > Filter

#Filter entries of age, only keep those records of which the values are >24

```
>>> df.filter(df["age"]>24).show()
```

## > Duplicate Values

```
>>> df = df.dropDuplicates()
```

## > Queries

```
>>> from pyspark.sql import functions as F
```

#### Select

```
>>> df.select("firstName").show() #Show all entries in firstName column
>>> df.select("firstName","lastName") \
    .show()
>>> df.select("firstName", #Show all entries in firstName, age and type
    "age",
    explode("phoneNumber") \
    .alias("contactInfo")) \
    .select("contactInfo.type",
    "firstName",
    "age") \
    .show()
>>> df.select(df["firstName"],df["age"]+ 1) #Show all entries in firstName and age,
    .show()                                add 1 to the entries of age
>>> df.select(df['age'] > 24).show() #Show all entries where age >24
```

#### When

```
>>> df.select("firstName", #Show firstName and 0 or 1 depending on age >30
    F.when(df.age > 30, 1) \
    .otherwise(0)) \
    .show()
>>> df[df.firstName.isin("Jane","Boris")] #Show firstName if in the given options
    .collect()
```

#### Like

```
>>> df.select("firstName", #Show firstName, and lastName is TRUE if lastName is like Smith
    df.lastName.like("Smith")) \
    .show()
```

#### Startswith - Endswith

```
>>> df.select("firstName", #Show firstName, and TRUE if lastName starts with Sm
    df.lastName \
    .startswith("Sm")) \
    .show()
>>> df.select(df.lastName.endswith("th")) #Show last names ending in th
    .show()
```

#### Substring

```
>>> df.select(df.firstName.substr(1, 3) \ #Return substrings of firstName
    .alias("name")) \
    .collect()
```

#### Between

```
>>> df.select(df.age.between(22, 24)) \ #Show age: values are TRUE if between 22 and 24
    .show()
```

## > Add, Update & Remove Columns

#### Adding Columns

```
>>> df = df.withColumn('city',df.address.city) \
    .withColumn('postalCode',df.address.postalCode) \
    .withColumn('state',df.address.state) \
    .withColumn('streetAddress',df.address.streetAddress) \
    .withColumn('telephoneNumber', explode(df.phoneNumber.number)) \
    .withColumn('telephoneType', explode(df.phoneNumber.type))
```

#### Updating Columns

```
>>> df = df.withColumnRenamed('telephoneNumber', 'phoneNumber')
```

#### Removing Columns

```
>>> df = df.drop("address", "phoneNumber")
>>> df = df.drop(df.address).drop(df.phoneNumber)
```

## > Missing & Replacing Values

```
>>> df.na.fill(50).show() #Replace null values
>>> df.na.drop().show() #Return new df omitting rows with null values
>>> df.na \ #Return new df replacing one value with another
    .replace(10, 20) \
    .show()
```

## > GroupBy

```
>>> df.groupBy("age")\ #Group by age, count the members in the groups
    .count() \
    .show()
```

## > Sort

```
>>> peopledf.sort(peopledf.age.desc()).collect()
>>> df.sort("age", ascending=False).collect()
>>> df.orderBy(["age", "city"], ascending=[0,1])\
    .collect()
```

## > Repartitioning

```
>>> df.repartition(10)\ #df with 10 partitions
    .rdd \
    .getNumPartitions()
>>> df.coalesce(1).rdd.getNumPartitions() #df with 1 partition
```

## > Running Queries Programmatically

#### Registering DataFrames as Views

```
>>> peopledf.createGlobalTempView("people")
>>> df.createTempView("customer")
>>> df.createOrReplaceTempView("customer")
```

#### Query Views

```
>>> df5 = spark.sql("SELECT * FROM customer").show()
>>> peopledf2 = spark.sql("SELECT * FROM global_temp.people")\
    .show()
```

## > Inspect Data

```
>>> df.dtypes #Return df column names and data types
```

```
>>> df.show() #Display the content of df
```

```
>>> df.head() #Return first n rows
```

```
>>> df.first() #Return first row
```

```
>>> df.take(2) #Return the first n rows >> df.schema Return the schema of df
```

```
>>> df.describe().show() #Compute summary statistics >> df.columns Return the columns of df
```

```
>>> df.count() #Count the number of rows in df
```

```
>>> df.distinct().count() #Count the number of distinct rows in df
```

```
>>> df.printSchema() #Print the schema of df
```

```
>>> df.explain() #Print the (logical and physical) plans
```

## > Output

#### Data Structures

```
>>> rdd1 = df.rdd #Convert df into an RDD
```

```
>>> df.toJSON().first() #Convert df into a RDD of string
```

```
>>> df.toPandas() #Return the contents of df as Pandas DataFrame
```

#### Write & Save to Files

```
>>> df.select("firstName", "city")\
```

```
    .write \
    .save("nameAndCity.parquet")
```

```
>>> df.select("firstName", "age")\
```

```
    .write \
    .save("namesAndAges.json", format="json")
```

## > Stopping SparkSession

```
>>> spark.stop()
```

| ALGORITHM             | DESCRIPTION                         | APPLICATIONS                                                                                                                                                                         | ADVANTAGES                                                                                                             | DISADVANTAGES                                                                                                                                               |
|-----------------------|-------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Linear Models         | <b>Linear Regression</b>            | A simple algorithm that models a linear relationship between inputs and a continuous numerical output variable                                                                       | <b>USE CASES</b><br>1. Stock price prediction<br>2. Predicting housing prices<br>3. Predicting customer lifetime value | 1. Explainable method<br>2. Interpretable results by its output coefficients<br>3. Faster to train than other machine learning models                       |
|                       | <b>Logistic Regression</b>          | A simple algorithm that models a linear relationship between inputs and a categorical output (1 or 0)                                                                                | <b>USE CASES</b><br>1. Credit risk score prediction<br>2. Customer churn prediction                                    | 1. Interpretable and explainable<br>2. Less prone to overfitting when using regularization<br>3. Applicable for multi-class predictions                     |
|                       | <b>Ridge Regression</b>             | Part of the regression family — it penalizes features that have low predictive outcomes by shrinking their coefficients closer to zero. Can be used for classification or regression | <b>USE CASES</b><br>1. Predictive maintenance for automobiles<br>2. Sales revenue prediction                           | 1. Less prone to overfitting<br>2. Best suited where data suffer from multicollinearity<br>3. Explainable & interpretable                                   |
|                       | <b>Lasso Regression</b>             | Part of the regression family — it penalizes features that have low predictive outcomes by shrinking their coefficients to zero. Can be used for classification or regression        | <b>USE CASES</b><br>1. Predicting housing prices<br>2. Predicting clinical outcomes based on health data               | 1. Less prone to overfitting<br>2. Can handle high-dimensional data<br>3. No need for feature selection                                                     |
| Supervised Learning   | <b>Decision Tree</b>                | Decision Tree models make decision rules on the features to produce predictions. It can be used for classification or regression                                                     | <b>USE CASES</b><br>1. Customer churn prediction<br>2. Credit score modeling<br>3. Disease prediction                  | 1. Explainable and interpretable<br>2. Can handle missing values                                                                                            |
|                       | <b>Random Forests</b>               | An ensemble learning method that combines the output of multiple decision trees                                                                                                      | <b>USE CASES</b><br>1. Credit score modeling<br>2. Predicting housing prices                                           | 1. Reduces overfitting<br>2. Higher accuracy compared to other models                                                                                       |
|                       | <b>Gradient Boosting Regression</b> | Gradient Boosting Regression employs boosting to make predictive models from an ensemble of weak predictive learners                                                                 | <b>USE CASES</b><br>1. Predicting car emissions<br>2. Predicting ride hailing fare amount                              | 1. Better accuracy compared to other regression models<br>2. It can handle multicollinearity<br>3. It can handle non-linear relationships                   |
|                       | <b>XGBoost</b>                      | Gradient Boosting algorithm that is efficient & flexible. Can be used for both classification and regression tasks                                                                   | <b>USE CASES</b><br>1. Churn prediction<br>2. Claims processing in insurance                                           | 1. Provides accurate results<br>2. Captures non linear relationships                                                                                        |
|                       | <b>LightGBM Regressor</b>           | A gradient boosting framework that is designed to be more efficient than other implementations                                                                                       | <b>USE CASES</b><br>1. Predicting flight time for airlines<br>2. Predicting cholesterol levels based on health data    | 1. Can handle large amounts of data<br>2. Computational efficient & fast training speed<br>3. Low memory usage                                              |
| Unsupervised Learning | <b>K-Means</b>                      | K-Means is the most widely used clustering approach—it determines K clusters based on euclidean distances                                                                            | <b>USE CASES</b><br>1. Customer segmentation<br>2. Recommendation systems                                              | 1. Scales to large datasets<br>2. Simple to implement and interpret<br>3. Results in tight clusters                                                         |
|                       | <b>Hierarchical Clustering</b>      | A "bottom-up" approach where each data point is treated as its own cluster—and then the closest two clusters are merged together iteratively                                         | <b>USE CASES</b><br>1. Fraud detection<br>2. Document clustering based on similarity                                   | 1. There is no need to specify the number of clusters<br>2. The resulting dendrogram is informative                                                         |
|                       | <b>Gaussian Mixture Models</b>      | A probabilistic model for modeling normally distributed clusters within a dataset                                                                                                    | <b>USE CASES</b><br>1. Customer segmentation<br>2. Recommendation systems                                              | 1. Computes a probability for an observation belonging to a cluster<br>2. Can identify overlapping clusters<br>3. More accurate results compared to K-means |
| Association           | <b>Apriori algorithm</b>            | Rule based approach that identifies the most frequent itemset in a given dataset where prior knowledge of frequent itemset properties is used                                        | <b>USE CASES</b><br>1. Product placements<br>2. Recommendation engines<br>3. Promotion optimization                    | 1. Results are intuitive and Interpretable<br>2. Exhaustive approach as it finds all rules based on the confidence and support                              |
|                       |                                     |                                                                                                                                                                                      |                                                                                                                        | 1. Generates many uninteresting itemsets<br>2. Computationally and memory intensive.<br>3. Results in many overlapping item sets                            |

# Python For Data Science Keras Cheat Sheet

Learn Keras online at [www.DataCamp.com](http://www.DataCamp.com)

## Keras

Keras is a powerful and easy-to-use deep learning library for Theano and TensorFlow that provides a high-level neural networks API to develop and evaluate deep learning models.

### A Basic Example

```
>>> import numpy as np
>>> from tensorflow.keras.models import Sequential
>>> from tensorflow.keras.layers import Dense
>>> data = np.random.random((1000,100))
>>> labels = np.random.randint(2,size=(1000,1))
>>> model = Sequential()
>>> model.add(Dense(32,
    activation='relu',
    input_dim=100))
>>> model.add(Dense(1, activation='sigmoid'))
>>> model.compile(optimizer='rmsprop',
    loss='binary_crossentropy',
    metrics=['accuracy'])
>>> model.fit(data,labels,epochs=10,batch_size=32)
>>> predictions = model.predict(data)
```

## Data

Your data needs to be stored as NumPy arrays or as a list of NumPy arrays. Ideally, you split the data in training and test sets, for which you can also resort to the `train_test_split` module of `sklearn.cross_validation`.

### Keras Data Sets

```
>>> from tensorflow.keras.datasets import boston_housing, mnist, cifar10, imdb
>>> (x_train,y_train),(x_test,y_test) = mnist.load_data()
>>> (x_train2,y_train2),(x_test2,y_test2) = boston_housing.load_data()
>>> (x_train3,y_train3),(x_test3,y_test3) = cifar10.load_data()
>>> (x_train4,y_train4),(x_test4,y_test4) = imdb.load_data(num_words=20000)
>>> num_classes = 10
```

### Other

```
>>> from urllib.request import urlopen
>>> data =
np.loadtxt(urlopen("http://archive.ics.uci.edu/ml/machine-learning-databases/pima-indians-diabetes/pima-indians-diabetes.data"),delimiter=",")
>>> X = data[:,0:8]
>>> y = data[:,8]
```

## Preprocessing

### Sequence Padding

```
>>> from tensorflow.keras.preprocessing import sequence
>>> x_train4 = sequence.pad_sequences(x_train4,maxlen=80)
>>> x_test4 = sequence.pad_sequences(x_test4,maxlen=80)
```

### One-Hot Encoding

```
>>> from tensorflow.keras.utils import to_categorical
>>> Y_train = to_categorical(y_train, num_classes)
>>> Y_test = to_categorical(y_test, num_classes)
>>> Y_train3 = to_categorical(y_train3, num_classes)
>>> Y_test3 = to_categorical(y_test3, num_classes)
```

## > Model Architecture

### Sequential Model

```
>>> from tensorflow.keras.models import Sequential
>>> model = Sequential()
>>> model2 = Sequential()
>>> model3 = Sequential()
```

### Multilayer Perceptron (MLP)

#### Binary Classification

```
>>> from tensorflow.keras.layers import Dense
>>> model.add(Dense(12,
    input_dim=8,
    kernel_initializer='uniform',
    activation='relu'))
>>> model.add(Dense(8,kernel_initializer='uniform',activation='relu'))
>>> model.add(Dense(1,kernel_initializer='uniform',activation='sigmoid'))
```

#### Multi-Class Classification

```
>>> from tensorflow.keras.layers import Dropout
>>> model.add(Dense(512,activation='relu',input_shape=(784,)))
>>> model.add(Dropout(0.2))
>>> model.add(Dense(512,activation='relu'))
>>> model.add(Dropout(0.2))
>>> model.add(Dense(10,activation='softmax'))
```

#### Regression

```
>>> model.add(Dense(64,activation='relu',input_dim=train_data.shape[1]))
>>> model.add(Dense(1))
```

### Convolutional Neural Network (CNN)

```
>>> from tensorflow.keras.layers import Activation,Conv2D,MaxPooling2D,Flatten
>>> model2.add(Conv2D(32,(3,3),padding='same',input_shape=x_train.shape[1:]))
>>> model2.add(Activation('relu'))
>>> model2.add(Conv2D(32,(3,3)))
>>> model2.add(Activation('relu'))
>>> model2.add(MaxPooling2D(pool_size=(2,2)))
>>> model2.add(Dropout(0.25))
>>> model2.add(Conv2D(64,(3,3), padding='same'))
>>> model2.add(Activation('relu'))
>>> model2.add(Conv2D(64,(3,3)))
>>> model2.add(Activation('relu'))
>>> model2.add(MaxPooling2D(pool_size=(2,2)))
>>> model2.add(Dropout(0.25))
>>> model2.add(Flatten())
>>> model2.add(Dense(512))
>>> model2.add(Activation('relu'))
>>> model2.add(Dropout(0.5))
>>> model2.add(Dense(num_classes))
>>> model2.add(Activation('softmax'))
```

### Recurrent Neural Network (RNN)

```
>>> from tensorflow.keras.layers import Embedding,LSTM
>>> model3.add(Embedding(20000,128))
>>> model3.add(LSTM(128,dropout=0.2,recurrent_dropout=0.2))
>>> model3.add(Dense(1,activation='sigmoid'))
```

## > Prediction

```
>>> model3.predict(x_test4, batch_size=32)
>>> model3.predict_classes(x_test4,batch_size=32)
```

Also see NumPy & Scikit-Learn

### Train and Test Sets

```
>>> from sklearn.model_selection import train_test_split
>>> X_train5,X_test5,y_train5,y_test5 = train_test_split(x, y,
    test_size=0.33,
    random_state=42)
```

### Standardization/Normalization

```
>>> from sklearn.preprocessing import StandardScaler
>>> scaler = StandardScaler().fit(x_train2)
>>> standardized_X = scaler.transform(x_train2)
>>> standardized_X_test = scaler.transform(x_test2)
```

## > Inspect Model

```
>>> model.output_shape #Model output shape
>>> model.summary() #Model summary representation
>>> model.get_config() #Model configuration
>>> model.get_weights() #List all weight tensors in the model
```

## > Compile Model

#### MLP: Binary Classification

```
>>> model.compile(optimizer='adam',
    loss='binary_crossentropy',
    metrics=['accuracy'])
```

#### MLP: Multi-Class Classification

```
>>> model.compile(optimizer='rmsprop',
    loss='categorical_crossentropy',
    metrics=['accuracy'])
```

#### MLP: Regression

```
>>> model.compile(optimizer='rmsprop',
    loss='mse',
    metrics=['mae'])
```

#### Recurrent Neural Network

```
>>> model3.compile(loss='binary_crossentropy',
    optimizer='adam',
    metrics=['accuracy'])
```

## > Model Training

```
>>> model3.fit(x_train4,
    y_train4,
    batch_size=32,
    epochs=15,
    verbose=1,
    validation_data=(x_test4,y_test4))
```

## > Evaluate Your Model's Performance

```
>>> score = model3.evaluate(x_test,
    y_test,
    batch_size=32)
```

## > Save/ Reload Models

```
>>> from keras.models import load_model
>>> model3.save('model_file.h5')
>>> my_model = load_model('my_model.h5')
```

## > Model Fine-tuning

### Optimization Parameters

```
>>> from tensorflow.keras.optimizers import RMSprop
>>> opt = RMSprop(lr=0.0001, decay=1e-6)
>>> model2.compile(loss='categorical_crossentropy',
    optimizer=opt,
    metrics=['accuracy'])
```

### Early Stopping

```
>>> from tensorflow.keras.callbacks import EarlyStopping
>>> early_stopping_monitor = EarlyStopping(patience=2)
>>> model3.fit(x_train4,
    y_train4,
    batch_size=32,
    epochs=15,
    validation_data=(x_test4,y_test4),
    callbacks=[early_stopping_monitor])
```